

A Proof of Lemma 1

Proof. It follows that

$$\begin{aligned}
\psi_{1,1}(\mathbf{z}_{1,1}) &\stackrel{(a)}{=} \psi_{1,0}(\mathbf{z}_{1,1}) + a_1(g(\mathbf{y}_{1,1}) + \langle \nabla g(\mathbf{y}_{1,1}), \mathbf{z}_{1,1} - \mathbf{y}_{1,1} \rangle + l(\mathbf{z}_{1,1})) \\
&\stackrel{(b)}{=} \frac{1}{2} \|\mathbf{z}_{1,1} - \mathbf{z}_{1,0}\|^2 + a_1(g(\mathbf{y}_{1,1}) + \langle \nabla g(\mathbf{y}_{1,1}), \mathbf{z}_{1,1} - \mathbf{y}_{1,1} \rangle + l(\mathbf{z}_{1,1})) \\
&\stackrel{(c)}{=} a_1 \left(g(\mathbf{y}_{1,1}) + \langle \nabla g(\mathbf{y}_{1,1}), \mathbf{z}_{1,1} - \mathbf{y}_{1,1} \rangle + \frac{1}{2a_1} \|\mathbf{z}_{1,1} - \mathbf{y}_{1,1}\|^2 + l(\mathbf{z}_{1,1}) \right) \\
&\stackrel{(d)}{=} a_1 \left(g(\mathbf{y}_{1,1}) + \langle \nabla g(\mathbf{y}_{1,1}), \mathbf{z}_{1,1} - \mathbf{y}_{1,1} \rangle + \frac{L}{2} \|\mathbf{z}_{1,1} - \mathbf{y}_{1,1}\|^2 + l(\mathbf{z}_{1,1}) \right) \\
&\stackrel{(e)}{\geq} a_1(g(\mathbf{z}_{1,1}) + l(\mathbf{z}_{1,1})) \\
&\stackrel{(f)}{=} A_1 f(\tilde{\mathbf{x}}_1),
\end{aligned}$$

where (a) is by definition of $\psi_{1,1}$, (b) is by the definition of $\psi_{1,0}$ and $\mathbf{z}_{1,0} = \tilde{\mathbf{x}}_0$, (c) is by the setting $\mathbf{y}_{1,1} = \mathbf{z}_{1,0}$ and simple rearrangement, (d) is by the setting $a_1 = \frac{1}{L}$, (e) is by Lemma 6, and (f) is by the setting $A_1 = a_1$ and $\tilde{\mathbf{x}}_1 = \mathbf{z}_{1,1}$. \square

B Proof of Lemma 2

Proof. As $l(\mathbf{z})$ is σ -strongly convex, by the definition of the sequence $\{\psi_{s,k}(\mathbf{z})\}$, $\psi_{s-1,m}(\mathbf{z})$ is $m + \sigma m \sum_{i=1}^{s-1} a_i = m(1 + \sigma A_{s-1})$ -strongly convex. Furthermore, we also know that $\psi_{s,k}(\mathbf{z}) (k \geq 0)$ is also at least $m(1 + \sigma A_{s-1})$ -strongly convex. So it follows that: $\forall k \geq 1$,

$$\begin{aligned}
\psi_{s,k}(\mathbf{z}_{s,k}) &\stackrel{(a)}{=} \psi_{s,k-1}(\mathbf{z}_{s,k}) + a_s(g(\mathbf{y}_{s,k}) + \langle \tilde{\nabla}_{s,k}, \mathbf{z}_{s,k} - \mathbf{y}_{s,k} \rangle + l(\mathbf{z}_{s,k})) \\
&\stackrel{(b)}{\geq} \psi_{s,k-1}(\mathbf{z}_{s,k-1}) + \frac{m(1 + \sigma A_{s-1})}{2} \|\mathbf{z}_{s,k} - \mathbf{z}_{s,k-1}\|^2 \\
&\quad + a_s(g(\mathbf{y}_{s,k}) + \langle \tilde{\nabla}_{s,k}, \mathbf{z}_{s,k} - \mathbf{y}_{s,k} \rangle + l(\mathbf{z}_{s,k})), \tag{19}
\end{aligned}$$

where (a) is by the definition of $\psi_{s,k}$ and (b) is by the optimality condition of $\mathbf{z}_{s,k-1}$ and the $m(1 + \sigma A_{s-1})$ -strong convexity of $\psi_{s,k-1}$. Then we have

$$\begin{aligned}
&a_s(g(\mathbf{y}_{s,k}) + \langle \tilde{\nabla}_{s,k}, \mathbf{z}_{s,k} - \mathbf{y}_{s,k} \rangle + l(\mathbf{z}_{s,k})) \\
&\stackrel{(a)}{=} a_s g(\mathbf{y}_{s,k}) + A_s \left\langle \tilde{\nabla}_{s,k}, \frac{a_s}{A_s} \mathbf{z}_{s,k} - \mathbf{y}_{s,k} + \frac{A_{s-1}}{A_s} \tilde{\mathbf{x}}_{s-1} \right\rangle \\
&\quad - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle + a_s l(\mathbf{z}_{s,k}) \\
&\stackrel{(b)}{\geq} a_s g(\mathbf{y}_{s,k}) + A_s \left\langle \tilde{\nabla}_{s,k}, \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \right\rangle - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle \\
&\quad + A_s l(\mathbf{y}_{s,k+1}) - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}), \tag{20}
\end{aligned}$$

where (a) is by the fact that $A_s = A_{s-1} + a_s$ and simple rearrangement and (b) is by $\mathbf{y}_{s,k+1} = \frac{A_{s-1}}{A_s} \tilde{\mathbf{x}}_{s-1} + \frac{a_s}{A_s} \mathbf{z}_{s,k}$ (which is by our definition of the sequence $\{\mathbf{y}_{s,k}\}$) and the convexity of $l(\mathbf{z})$.

Meanwhile, by our setting in Step 5 of Algorithm 1, $A_s = A_{s-1} + \sqrt{\frac{mA_{s-1}(1+\sigma A_{s-1})}{2L}}$ and also $a_s = A_s - A_{s-1}$, we have

$$\frac{mA_s(1 + \sigma A_{s-1})}{a_s^2} = \frac{2A_s}{A_{s-1}} L \geq \left(1 + \frac{A_s}{A_{s-1}}\right) L. \tag{21}$$

Then by combining (19) and (20), it follows that

$$\begin{aligned}
& \psi_{s,k}(\mathbf{z}_{s,k}) - \psi_{s,k-1}(\mathbf{z}_{s,k-1}) \\
\geq & a_s g(\mathbf{y}_{s,k}) + A_s \langle \tilde{\nabla}_{s,k}, \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle \\
& + A_s l(\mathbf{y}_{s,k+1}) - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}) + \frac{m(1 + \sigma A_{s-1})}{2} \|\mathbf{z}_{s,k} - \mathbf{z}_{s,k-1}\|^2 \\
\stackrel{(a)}{=} & a_s g(\mathbf{y}_{s,k}) + A_s \langle \tilde{\nabla}_{s,k}, \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle \\
& + A_s l(\mathbf{y}_{s,k+1}) - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}) + \frac{mA_s^2(1 + \sigma A_{s-1})}{2a_s^2} \|\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k}\|^2 \\
\stackrel{(b)}{\geq} & a_s g(\mathbf{y}_{s,k}) + A_s \left(\langle \tilde{\nabla}_{s,k}, \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle \right. \\
& \quad \left. + \left(1 + \frac{A_s}{A_{s-1}}\right) \frac{L}{2} \|\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k}\|^2 + l(\mathbf{y}_{s,k+1}) \right) \\
& - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}),
\end{aligned}$$

where (a) is by the fact $\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} = \frac{a_s}{A_s}(\mathbf{z}_{s,k} - \mathbf{z}_{s,k-1})$ and (b) is by (21). Then we have

$$\begin{aligned}
& \langle \tilde{\nabla}_{s,k}, \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle + \left(1 + \frac{A_s}{A_{s-1}}\right) \frac{L}{2} \|\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k}\|^2 + l(\mathbf{y}_{s,k+1}) \\
= & \langle \nabla g(\mathbf{y}_{s,k}), \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle + \frac{L}{2} \|\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k}\|^2 + l(\mathbf{y}_{s,k+1}) \\
& + \langle \tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k}), \mathbf{y}_{s,k+1} - \mathbf{y}_{s,k} \rangle + \frac{A_s L}{2A_{s-1}} \|\mathbf{y}_{s,k+1} - \mathbf{y}_{s,k}\|^2 \\
\stackrel{(a)}{\geq} & g(\mathbf{y}_{s,k+1}) - g(\mathbf{y}_{s,k}) + l(\mathbf{y}_{s,k+1}) - \frac{A_{s-1}}{2A_s L} \|\tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k})\|^2 \\
= & f(\mathbf{y}_{s,k+1}) - g(\mathbf{y}_{s,k}) - \frac{A_{s-1}}{2A_s L} \|\tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k})\|^2 \tag{22}
\end{aligned}$$

where (a) is by Lemma 6 and the Young's inequality $\langle \mathbf{a}, \mathbf{b} \rangle \geq -\frac{1}{2}\|\mathbf{a}\|^2 - \frac{1}{2}\|\mathbf{b}\|^2$. So we have

$$\begin{aligned}
& \psi_{s,k}(\mathbf{z}_{s,k}) - \psi_{s,k-1}(\mathbf{z}_{s,k-1}) \\
\geq & a_s g(\mathbf{y}_{s,k}) + A_s \left(f(\mathbf{y}_{s,k+1}) - g(\mathbf{y}_{s,k}) - \frac{A_{s-1}}{2A_s L} \|\tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k})\|^2 \right) \\
& - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}). \tag{23}
\end{aligned}$$

□

C Proof of Lemma 3

Proof. Taking expectation on the randomness over the choice of i , we have

$$\begin{aligned}
\mathbb{E}[\|\tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k})\|^2] &= \mathbb{E}[\|\nabla g_i(\mathbf{y}_{s,k}) - \nabla g_i(\tilde{\mathbf{x}}_{s-1}) + \boldsymbol{\mu}_{s-1} - \nabla g(\mathbf{y}_{s,k})\|^2] \\
&= \mathbb{E}[\|\nabla g_i(\mathbf{y}_{s,k}) - \nabla g_i(\tilde{\mathbf{x}}_{s-1}) + \nabla g(\tilde{\mathbf{x}}_{s-1}) - \nabla g(\mathbf{y}_{s,k})\|^2] \\
&= \mathbb{E}[\|\nabla g_i(\mathbf{y}_{s,k}) - \nabla g_i(\tilde{\mathbf{x}}_{s-1})\|^2] - \|\nabla g(\tilde{\mathbf{x}}_{s-1}) - \nabla g(\mathbf{y}_{s,k})\|^2 \\
&\leq \mathbb{E}[\|\nabla g_i(\mathbf{y}_{s,k}) - \nabla g_i(\tilde{\mathbf{x}}_{s-1})\|^2] \\
&\stackrel{(a)}{\leq} \mathbb{E}[2L(g_i(\tilde{\mathbf{x}}_{s-1}) - g_i(\mathbf{y}_{s,k}) - \langle \nabla g_i(\mathbf{y}_{s,k}), \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle)] \\
&= 2L(g(\tilde{\mathbf{x}}_{s-1}) - g(\mathbf{y}_{s,k}) - \langle \nabla g(\mathbf{y}_{s,k}), \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle),
\end{aligned}$$

where (a) is by Lemma 6. □

D Proof of Lemma 4

Proof. By Lemma 2 and taking expectation on the randomness over the choice of i , we have

$$\begin{aligned}
& \mathbb{E}[\psi_{s,k}(\mathbf{z}_{s,k}) - \psi_{s,k-1}(\mathbf{z}_{s,k-1})] \\
\geq & \mathbb{E}\left[a_s g(\mathbf{y}_{s,k}) + A_s \left(f(\mathbf{y}_{s,k+1}) - g(\mathbf{y}_{s,k}) - \frac{A_{s-1}}{2A_s L} \|\tilde{\nabla}_{s,k} - \nabla g(\mathbf{y}_{s,k})\|^2 \right) \right. \\
& \left. - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}) \right] \\
\stackrel{(a)}{\geq} & \mathbb{E}\left[a_s g(\mathbf{y}_{s,k}) \right. \\
& \left. + A_s (f(\mathbf{y}_{s,k+1}) - g(\mathbf{y}_{s,k})) - A_{s-1} (g(\tilde{\mathbf{x}}_{s-1}) - g(\mathbf{y}_{s,k}) - \langle \nabla g(\mathbf{y}_{s,k}), \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle) \right. \\
& \left. - A_{s-1} \langle \tilde{\nabla}_{s,k}, \tilde{\mathbf{x}}_{s-1} - \mathbf{y}_{s,k} \rangle - A_{s-1} l(\tilde{\mathbf{x}}_{s-1}) \right] \\
\stackrel{(b)}{=} & \mathbb{E}[A_s f(\mathbf{y}_{s,k+1})] - A_{s-1} f(\tilde{\mathbf{x}}_{s-1}), \tag{24}
\end{aligned}$$

where (a) is by Lemma 3, and (b) is by $\mathbb{E}[\tilde{\nabla}_{s,k}] = \nabla g(\mathbf{y}_{s,k})$, $A_s = A_{s-1} + a_s$ and $f(\mathbf{x}) = g(\mathbf{x}) + l(\mathbf{x})$.

Summing (24) from $k = 1$ to m , by the setting for $s \geq 2$, $\psi_{s+1,0} := \psi_{s,m}$ and $\mathbf{z}_{s+1,0} := \mathbf{z}_{s,m}$, we have

$$\begin{aligned}
\mathbb{E}[\psi_{s+1,0}(\mathbf{z}_{s+1,0}) - \psi_{s,0}(\mathbf{z}_{s,0})] &= \mathbb{E}[\psi_{s,m}(\mathbf{z}_{s,m}) - \psi_{s,0}(\mathbf{z}_{s,0})] \\
&\geq \mathbb{E}\left[A_s \sum_{k=1}^m f(\mathbf{y}_{s,k+1}) - mA_{s-1} f(\tilde{\mathbf{x}}_{s-1}) \right] \\
&\stackrel{(a)}{\geq} \mathbb{E}\left[mA_s f(\tilde{\mathbf{x}}_s) - mA_{s-1} f(\tilde{\mathbf{x}}_{s-1}) \right], \tag{25}
\end{aligned}$$

where (a) is by the convexity of $f(\mathbf{z})$ and the fact of $\tilde{\mathbf{x}}_s = \frac{1}{m} \sum_{k=1}^m \mathbf{y}_{s,k+1}$ (which is in turn by the definition of $\tilde{\mathbf{x}}_s = \frac{A_{s-1}}{A_s} \tilde{\mathbf{x}}_{s-1} + \frac{a_s}{mA_s} \sum_{k=1}^m \mathbf{z}_{s,k}$ and the definition of $\mathbf{y}_{s,k}$)

□

E Proof of Lemma 5

Proof. $\forall s \geq 2$, taking expectation on the choice of i in the k -th iteration of the s -th epoch, we have $\forall \mathbf{z}$,

$$\begin{aligned}
\mathbb{E}[\psi_{s,k}(\mathbf{z})] &= \mathbb{E}[\psi_{s,k-1}(\mathbf{z}) + a_s (g(\mathbf{y}_{s,k}) + \langle \tilde{\nabla}_{s,k}, \mathbf{z} - \mathbf{y}_{s,k} \rangle + l(\mathbf{z}))] \\
&\stackrel{(a)}{=} \psi_{s,k-1}(\mathbf{z}) + a_s (g(\mathbf{y}_{s,k}) + \langle \nabla g(\mathbf{y}_{s,k}), \mathbf{z} - \mathbf{y}_{s,k} \rangle + l(\mathbf{z})) \\
&\stackrel{(b)}{\leq} \psi_{s,k-1}(\mathbf{z}) + a_s (g(\mathbf{z}) + l(\mathbf{z})) \\
&= \psi_{s,k-1}(\mathbf{z}) + a_s f(\mathbf{z}), \tag{26}
\end{aligned}$$

where (a) is by the fact $\mathbb{E}[\tilde{\nabla}_{s,k}] = \nabla g(\mathbf{y}_{s,k})$, and (b) is by the convexity of $g(\mathbf{z})$. Then taking expectation from the randomness of the epoch s and telescoping (26) from $k = 1$ to m , we have

$$\begin{aligned}
\mathbb{E}[\psi_{s,m}(\mathbf{z})] &\leq \psi_{s,0}(\mathbf{z}) + ma_s f(\mathbf{z}) \\
&= \begin{cases} \psi_{s-1,m}(\mathbf{z}) + ma_s f(\mathbf{z}), & s \geq 3 \\ m\psi_{1,1}(\mathbf{z}) + ma_2 f(\mathbf{z}), & s = 2. \end{cases} \tag{27}
\end{aligned}$$

Then taking expectation from the randomness of all the history from $i = 3$ and telescoping (27) to some $s \geq 3$, we have

$$\mathbb{E}[\psi_{s,m}(\mathbf{z})] \leq \psi_{2,m}(\mathbf{z}) + m \sum_{i=3}^s a_i f(\mathbf{z}). \tag{28}$$

Meanwhile taking expectation from the randomness of epoch $s = 2$, we have

$$\begin{aligned}
\mathbb{E}[\psi_{2,m}(\mathbf{z})] &\leq m\psi_{1,1}(\mathbf{z}) + ma_2f(\mathbf{z}) \\
&= m(\psi_{1,0}(\mathbf{z}) + a_1(g(\mathbf{y}_{1,1}) + \langle \nabla g(\mathbf{y}_{1,1}), \mathbf{z} - \mathbf{y}_{1,1} \rangle + l(\mathbf{z}))) + ma_2f(\mathbf{z}) \\
&\stackrel{(a)}{\leq} m\left(\frac{1}{2}\|\mathbf{z} - \tilde{\mathbf{x}}_0\|^2 + a_1(g(\mathbf{z}) + l(\mathbf{z}))\right) + ma_2f(\mathbf{z}) \\
&= m(a_1 + a_2)f(\mathbf{z}) + \frac{m}{2}\|\mathbf{z} - \tilde{\mathbf{x}}_0\|^2,
\end{aligned} \tag{29}$$

where (a) is by the convexity of $g(\mathbf{z})$ and $\psi_{1,0}(\mathbf{z}) = \frac{1}{2}\|\mathbf{z} - \tilde{\mathbf{x}}_0\|^2$.

So combining (28) and (29), we have: $\forall s \geq 2$,

$$\begin{aligned}
\mathbb{E}[\psi_{s,m}(\mathbf{z})] &\leq m \sum_{i=1}^s a_s f(\mathbf{z}) + \frac{m}{2}\|\mathbf{z} - \tilde{\mathbf{x}}_0\|^2 \\
&\stackrel{(a)}{=} mA_s f(\mathbf{z}) + \frac{m}{2}\|\mathbf{z} - \tilde{\mathbf{x}}_0\|^2,
\end{aligned} \tag{30}$$

where (a) is by our setting $a_s = A_s - A_{s-1}$ and $A_0 = 0$.

Then by (30) and the optimality of $\mathbf{z}_{s,m}$, we have $\psi_{s,m}(\mathbf{z}_{s,m}) \leq \psi_{s,m}(\mathbf{x}^*)$ and thus

$$\mathbb{E}[\psi_{s,m}(\mathbf{z}_{s,m})] \leq \psi_{s,m}(\mathbf{x}^*) \leq mA_s f(\mathbf{x}^*) + \frac{m}{2}\|\mathbf{x}^* - \tilde{\mathbf{x}}_0\|^2. \tag{31}$$

□

F The Lower Bounds for the A_s in Theorem 1

Proof. In the following, we give the lower bound of A_s by the condition in Step 6 of Algorithm 1 and $A_1 = a_1 = \frac{1}{L}$. To show the lower bound by the first term in (10), we know that

$$A_s = A_{s-1} + \sqrt{\frac{mA_{s-1}(1 + \sigma A_{s-1})}{2L}} \geq \sqrt{\frac{mA_{s-1}(1 + \sigma A_{s-1})}{2L}} \geq \sqrt{\frac{mA_{s-1}}{2L}}, \tag{32}$$

so we have

$$\frac{2LA_s}{m} \geq \left(\frac{2LA_{s-1}}{m}\right)^{\frac{1}{2}} \geq \left(\frac{2LA_1}{m}\right)^{2^{-(s-1)}}. \tag{33}$$

Then by the setting $A_1 = \frac{1}{L}$, we have

$$A_s \geq \frac{m}{2L} \left(\frac{2}{m}\right)^{2^{-(s-1)}}. \tag{34}$$

Meanwhile, for $s \geq 2$, we also have

$$\begin{aligned}
A_s &\geq A_{s-1} + \sqrt{\frac{mA_{s-1}(1 + \sigma A_{s-1})}{2L}} \geq A_{s-1} + \sqrt{\frac{m\sigma}{2L}} A_{s-1} = \left(1 + \sqrt{\frac{m\sigma}{2L}}\right) A_{s-1} \\
&\geq \left(1 + \sqrt{\frac{m\sigma}{2L}}\right)^{s-1} A_1 \\
&= \frac{1}{L} \left(1 + \sqrt{\frac{m\sigma}{2L}}\right)^{s-1}.
\end{aligned} \tag{35}$$

Thus the lower bounds in (10) are proved.

Then with $s_0 = 1 + \lceil \log_2 \log_2(m/2) \rceil$, we have

$$\begin{aligned}
A_{s_0} &\geq \frac{m}{2L} \left(\frac{2}{m}\right)^{2^{-(s_0-1)}} \geq \frac{m}{2L} \left(\frac{2}{m}\right)^{2^{-\lceil \log_2 \log_2(m/2) \rceil}} \geq \frac{m}{2L} \left(\frac{2}{m}\right)^{2^{-\log_2 \log_2(m/2)}} \\
&= \frac{m}{4L}.
\end{aligned}$$

Meanwhile for $s \geq s_0 + 1$, we have

$$A_s \geq A_{s-1} + \sqrt{\frac{mA_{s-1}(1 + \sigma A_{s-1})}{2L}} \geq A_{s-1} + \sqrt{\frac{mA_{s-1}}{2L}}. \quad (36)$$

Thus we can use the mathematical induction method to prove the first lower bound in (11): $\forall s \geq s_0, A_s \geq \frac{m}{32L} (s - s_0 + 2\sqrt{2})^2$.

Firstly, for $s = s_0$, we have $A_s \geq \frac{m}{4L} = \frac{m}{32L} (2\sqrt{2})^2$.

Then assume that for an $s \geq s_0 + 1$, $A_{s-1} \geq \frac{m}{32L} (s - 1 - s_0 + 2\sqrt{2})^2$, then

$$\begin{aligned} A_s &\geq A_{s-1} + \sqrt{\frac{mA_{s-1}}{2L}} \geq \frac{m}{32L} (s - s_0 + 2\sqrt{2})^2 + \frac{m}{16L} (s - s_0) + \frac{m}{32L} (4\sqrt{2} - 3) \\ &\geq \frac{m}{32L} (s - s_0 + 2\sqrt{2})^2. \end{aligned} \quad (37)$$

Thus the first lower bound in (11) is proved.

Meanwhile, for $s \geq s_0 + 1$, we also have

$$\begin{aligned} A_s &\geq A_{s-1} + \sqrt{\frac{mA_{s-1}(1 + \sigma A_{s-1})}{2L}} \geq A_{s-1} + \sqrt{\frac{m\sigma}{2L}} A_{s-1} = \left(1 + \sqrt{\frac{m\sigma}{2L}}\right) A_{s-1} \\ &\geq \left(1 + \sqrt{\frac{m\sigma}{2L}}\right)^{s-s_0} A_{s_0} \\ &\geq \frac{m}{4L} \left(1 + \sqrt{\frac{m\sigma}{2L}}\right)^{s-s_0}. \end{aligned} \quad (38)$$

Thus the second lower bound in (11) is proved. □

G An Auxiliary Lemma

By Assumption 1 and [25], we have Lemma 6.

Lemma 6. Under Assumption 1, $\forall \mathbf{x}, \mathbf{y}$,

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (39)$$

and $\forall i \in [n], \forall \mathbf{x}, \mathbf{y}$,

$$\|\nabla g_i(\mathbf{y}) - \nabla g_i(\mathbf{x})\|^2 \leq 2L(g_i(\mathbf{y}) - g_i(\mathbf{x}) - \langle \nabla g_i(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle). \quad (40)$$

Under Assumption 1, Lemma 6 are classical results in convex optimization. For completeness, we provide the proof of Lemma 6 here.

Proof of Lemma 6. By Assumption 1, $\forall i \in [n], g_i(\mathbf{x})$ satisfies $\forall \mathbf{x}, \mathbf{y}, \|\nabla g_i(\mathbf{x}) - \nabla g_i(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$. As a result, we have

$$\begin{aligned} \|\nabla g(\mathbf{x}) - \nabla g(\mathbf{y})\| &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{x}) - \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{y}) \right\| \\ &\leq \frac{1}{n} \sum_{i=1}^n \|\nabla g_i(\mathbf{x}) - \nabla g_i(\mathbf{y})\| \\ &\leq L\|\mathbf{x} - \mathbf{y}\|. \end{aligned} \quad (41)$$

The we have

$$\begin{aligned}
g(\mathbf{y}) &= g(\mathbf{x}) + \int_0^1 \langle \nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle d\tau \\
&= g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \langle \nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau. \quad (42)
\end{aligned}$$

Then it follow that

$$\begin{aligned}
g(\mathbf{y}) - g(\mathbf{x}) - \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle &\leq \left| \int_0^1 \langle \nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau \right| \\
&\leq \int_0^1 |\langle \nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| d\tau \\
&\leq \int_0^1 \|\nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| d\tau \\
&\leq \int_0^1 L\tau \|\mathbf{y} - \mathbf{x}\|^2 d\tau \\
&= \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (43)
\end{aligned}$$

Thus we obtain (39).

Then denote $\forall i \in [n]$, $\phi_i(\mathbf{y}) = g_i(\mathbf{y}) - g_i(\mathbf{x}) - \langle \nabla g_i(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$. Obviously $\phi_i(\mathbf{y})$ is also L -smooth. One can check that $\nabla g_i(\mathbf{x}) = 0$ and so that $\min_{\mathbf{y}} \phi_i(\mathbf{y}) = \phi_i(\mathbf{x}) = 0$, which implies that

$$\begin{aligned}
\phi_i(\mathbf{x}) &\leq \phi_i\left(\mathbf{y} - \frac{1}{L} \nabla \phi_i(\mathbf{y})\right) \\
&= \phi_i(\mathbf{y}) + \int_0^1 \left\langle \nabla \phi_i\left(\mathbf{y} - \frac{\tau}{L} \nabla \phi_i(\mathbf{y})\right), -\frac{1}{L} \nabla \phi_i(\mathbf{y}) \right\rangle d\tau \\
&= \phi_i(\mathbf{y}) + \left\langle \nabla \phi_i(\mathbf{y}), -\frac{1}{L} \nabla \phi_i(\mathbf{y}) \right\rangle + \int_0^1 \left\langle \nabla \phi_i\left(\mathbf{y} - \frac{\tau}{L} \nabla \phi_i(\mathbf{y})\right) - \nabla \phi_i(\mathbf{y}), -\frac{1}{L} \nabla \phi_i(\mathbf{y}) \right\rangle d\tau \\
&\leq \phi_i(\mathbf{y}) - \frac{1}{L} \|\nabla \phi_i(\mathbf{y})\|^2 + \int_0^1 L \left\| \frac{\tau}{L} \nabla \phi_i(\mathbf{y}) \right\| \left\| \frac{1}{L} \nabla \phi_i(\mathbf{y}) \right\| d\tau \\
&\leq \phi_i(\mathbf{y}) - \frac{1}{2L} \|\nabla \phi_i(\mathbf{y})\|^2. \quad (44)
\end{aligned}$$

Then we have $\|\nabla \phi_i(\mathbf{y})\|^2 \leq 2L(\phi_i(\mathbf{y}) - \phi_i(\mathbf{x}))$. Then by the definition of $\phi_i(\mathbf{y})$, we obtain (40). \square

H Experimental Details and Supplementary Experiments

Besides running binary classification experiments on the two datasets a9a and covtype, we also run multi-class classification experiments on mnist and cifar10. The problem we solve is the ℓ_2 -norm regularized (multinomial) logistic regression problem:

$$\min_{\mathbf{w} \in \mathbb{R}^{d \times (c-1)}} f(\mathbf{w}) := \frac{1}{n} \sum_{j=1}^n \left(-\sum_{i=1}^{c-1} y_j^{(i)} \mathbf{w}^{(i)T} \mathbf{x}_j + \log \left(1 + \sum_{i=1}^{c-1} \exp(\mathbf{w}^{(i)T} \mathbf{x}_j) \right) \right) + \frac{\lambda}{2} \sum_{i=1}^{c-1} \|\mathbf{w}^{(i)}\|_2^2, \quad (45)$$

where n is the number of samples, $c \in \{2, 3, \dots\}$ denotes the number of class (for a9a and covtype, $c = 2$; for mnist and cifar10, $c = 10$), $\lambda \geq 0$ denotes the regularization parameter, $\mathbf{y}_j = (y_j^{(1)}, y_j^{(2)}, \dots, y_j^{(c-1)})^T$ is a one-hot vector or zero vector¹¹, and $\mathbf{w} := (\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(c-1)}) \in \mathbb{R}^{d \times (c-1)}$ denotes the variable to optimize. For the two-class datasets ‘‘a9a’’ and ‘‘covtype’’, we have presented our results by choosing the regularization parameter $\lambda \in \{0, 10^{-8}, 10^{-4}\}$. For the ten-class datasets ‘‘mnist’’ and ‘‘cifar10’’, we choose $\lambda \in \{0, 10^{-6}, 10^{-3}\}$.

¹¹Zero vector denotes the class of the j -th sample is c .

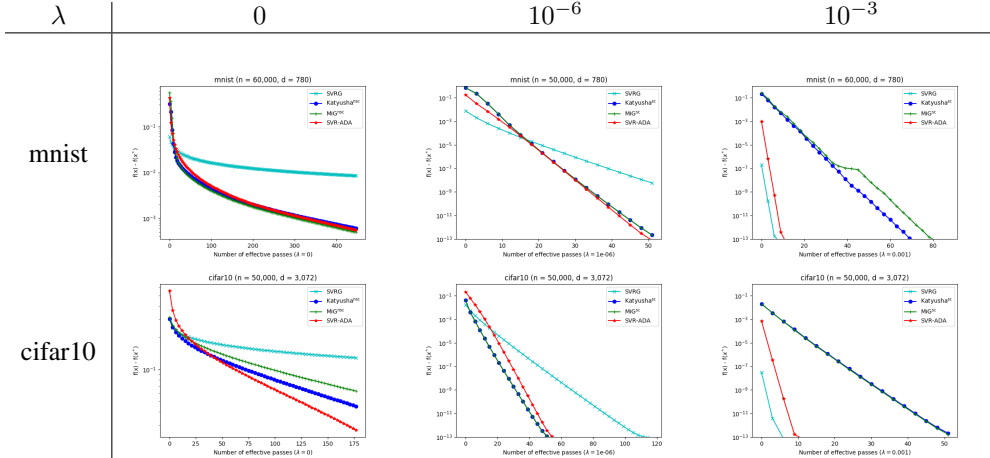


Figure 2: Comparing VRADA with SVRG, Katyusha and MiG on ℓ_2 -norm regularized multinomial logistic regression problems. The horizontal axis is the number of passes through the entire dataset, and the vertical axis is the optimality gap $f(\mathbf{x}) - f(\mathbf{x}^*)$.

For the four algorithms we compare, the common parameter to tune is the parameter *w.r.t.* Lipschitz constant¹², which is tuned in $\{0.0125, 0.025, 0.05, 0.1, 0.25, 0.5\}$.¹³ All four algorithms are implemented in C++ under the same framework, while the figures are produced using Python.

As we see, despite there are some minor differences among different tasks/datasets shown in Figure 1 and Figure 2, the general behaviors are still very consistent. From both figures, our method VRADA is competitive with other two accelerated methods, and is much faster than the non-accelerated SVRG algorithm in the general convex setting and the strongly convex setting with a large conditional number. Meanwhile, in the strongly convex setting with a small conditional number, VRADA is still competitive with the non-accelerated SVRG algorithm and much faster than the other two accelerated algorithms of Katyusha^{sc} and MiG^{sc}.

¹²For logistic regression with normalized data, the Lipschitz constant is globally upper bounded [39] by 1/4, but in practice we can use a smaller one than 1/4.

¹³In our experiments, due to the normalization of datasets, all the four algorithms will diverge when the parameter is less than 0.0125. Otherwise, they always converge if the parameter is less than 0.5.